

ITERATED GEOMETRIC HARMONICS FOR DATA IMPUTATION AND RECONSTRUCTION OF MISSING DATA

CHAD ECKMAN, JONATHAN A. LINDGREN, ERIN P. J. PEARSE, DAVID J. SACCO,
ZACHARIAH ZHANG

ABSTRACT. The method of geometric harmonics is adapted to the situation of incomplete data by means of the iterated geometric harmonics (IGH) scheme. The method is tested on natural and synthetic data sets with 50–500 data points and dimensionality of 400–10,000. Experiments suggest that the algorithm converges to a near optimal solution within 4–6 iterations, at runtimes of less than 30 minutes on a medium-grade desktop computer. The imputation of missing data values is applied to collections of damaged images (suffering from data annihilation rates of up to 70%) which are reconstructed with a surprising degree of accuracy.

1. INTRODUCTION

The method of geometric harmonics was introduced in the thesis of S. Lafon [Laf] (see also [CL2]) as a method of extending empirical functions defined on a dataset, considered as a point cloud $X = \{\mathbf{x}_i\}_{i=1}^n \subseteq \mathbb{R}^d$, where d is the dimension (number of parameters or characteristics) observed for each data point (observation) \mathbf{x}_i . In other words, consider a subset $A \subseteq \{1, \dots, n\}$ and let $X_A = \{\mathbf{x}_i \in X : i \in A\}$. Then given a function $f : X_A \rightarrow \mathbb{R}$, the goal is to construct an extension $\hat{f} : X \rightarrow \mathbb{R}$ for forecasting, classification, or other machine learning or statistical analysis purposes.

In the present paper, we do not consider inferring the value of some extrinsic function defined on the data (as in a classification or clustering problem) but turn the mechanism on the dataset itself. I.e., if one considers f as one column of the dataset, the task of extending the partially defined function $f : X_A \rightarrow \mathbb{R}$ amounts to imputing the missing values in a dataset where missing values occur only in that one column. In this paper, we introduce the method of iterated geometric harmonics (IGH) which uses an iteration scheme to adapt the method of geometric harmonics to more general situations of incomplete data (i.e., missing values in many/all columns). The method allows for reconstruction of datasets which are incomplete due to the presence of missing values (due to recording error, transmission loss, etc.). The datasets should have the property that they can be thought of as samples from some underlying topological manifold; i.e., the data should be comprised of one or more clusters of points which are “close” in the general sense of manifold learning theory. We focus on *missing* data, not *noisy* data; i.e., we assume that it is known which data need to be reconstructed.

Geometric harmonics is parameter-free in the sense that no regression model is assumed: the technique constructs the extension directly from the geometry of the dataset. Nonetheless, the analyst still must choose some way to measure similarities between data points (observations) and this amounts to determining a particular kernel for an integral operator, i.e., a positive semidefinite symmetric function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$. Some care must be taken when choosing the kernel; we give some suggestions for kernel selection below. We test our approach on natural and synthetic data sets and some conditions for convergence of

Key words and phrases. Data reconstruction, missing data, data imputation, geometric harmonics, diffusion map, machine learning, graph-theoretic models, graph algorithm, inference model, image processing.

The authors are grateful to support from the College of Science and Mathematics Summer Research Program at Cal Poly, San Luis Obispo.

the iteration are discussed. In most cases we find the iteration converges remarkably fast; typically no more than about 5 iterates are required.

1.1. Applications of IGH. IGH was developed for imputing missing data prior to scientific/statistical analysis; see the example of weather data given below. Conventional statistical software do not cope well with missing data; typically the analyst is required to invoke some imputation procedure or discard incomplete data points. Discarding incomplete data is a poor choice, as it can easily bias the remaining data or leave the analyst with too little data for a proper analysis; see [End, LR]. The present state-of-the-art technique for dealing with missing data is Multiple Imputation, in which multiple complete versions of the data are simulated by filling in missing entries in a purely stochastic manner; the analyses of these different simulated versions are then averaged (pooled). Unfortunately, this is essentially a linear technique and does not perform well when the data has a strongly nonlinear structure. See [End, LR, vB] for details.

It will be clear from the examples below that additional applications of IGH include image processing, especially video reconstruction, as video media naturally contains large numbers of similar images (frames). Some samples of reconstructed video appear online at

<http://www.calpoly.edu/~epearse/video.html>.

Based on the given examples, IGH clearly has potential applications for security, law enforcement, and the military, as well as reconstruction of archival footage and other tasks.

2. ITERATED GEOMETRIC HARMONICS

2.1. Geometric harmonics. The term “geometric harmonics” refers to an eigenbasis for an incomplete Gram matrix, constructed empirically by means of a variant of the Nyström method; cf. [CL2]. The Nyström method is a quadrature rule that effectively allows one to approximate the solution to an integral equation by subsampling: the integral is replaced by a sum of function values (the function is evaluated on some sample points), each one of which is weighted according to the quadrature rule. In the classical case, the subsample may be selected manually for maximal effectiveness; in modern applications (esp. to machine learning problems, note that the Nyström formula has been shown to be equivalent to kernel PCA projection [SSSM]), the subsample may be selected randomly as part of a Monte Carlo scheme. For examples related to spectral partitioning problems in image segmentation, see [FBCM, FBM]; for examples related to the extension of empirical eigenfunctions to data outside an original sample, see [BN, BNS, BPV⁺, Laf]; for examples related to the efficiency of kernel-based machine learning methods, see [WS, DM, GM].

In the context of geometric harmonics, one does not get to choose the subsample: it consists of those rows A of the data set for which the partially defined function $f : X_A \rightarrow \mathbb{R}$ is defined. For a dataset with n records, the method of geometric harmonics requires computing an $n \times n$ Gram matrix K with entries

$$K(i, j) = k(\mathbf{x}_i, \mathbf{x}_j), \quad (1)$$

where $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ is some symmetric, nonnegative, and typically positive semidefinite function. For example, for homogeneous data sampled uniformly from a submanifold of \mathbb{R}^d , we have used the Gaussian kernel

$$k(\mathbf{x}, \mathbf{y}) = \frac{1}{\sigma\sqrt{2\pi}} \exp(-\tfrac{1}{2}\|\mathbf{x} - \mathbf{y}\|_2^2/\sigma^2).$$

While the positive semidefinite property is not strictly necessary, it puts the method on firm theoretical footing within the context of reproducing kernel Hilbert spaces; cf. [Aro].

Geometric harmonics applies the Nyström method to the eigenequations

$$\sum_j K(i, j)\psi(j) = \lambda\psi(i). \quad (2)$$

The eventual application of the eigenbasis is to the extension of a function $f : X_A \rightarrow \mathbb{R}$ to a function $\hat{f} : X \rightarrow \mathbb{R}$, in which case the operator $K : f \mapsto Kf$ is not a square matrix. Therefore, one must consider diagonalizing the $|A| \times |A|$ matrix K^*K , and it can be shown that the adjoint K^* is simply the restriction operator; cf. [CL2, Lem. 1] or [Laf, Lem. 13]. This leads to solving (2) where j ranges only over A . However, it is observed in [Laf, §3.2] that the left side of (2) makes sense for any $i \in \{1, \dots, n\}$; only j needs to be restricted to A . Thus, one can turn the equation around and *define* the geometric harmonics Ψ_ℓ , for $\ell = 1, \dots, |A|$ by

$$\Psi_\ell(i) := \frac{1}{\lambda_\ell} \sum_{j \in A} K(i, j) \psi_\ell(j). \quad (3)$$

The key point is that (3) defines $\Psi_\ell(i)$ for *all* $i = 1, \dots, n$. Formula (3) generates a family of functions Ψ_ℓ which are maximally concentrated on A in a sense that generalizes the properties of the prolate spheroidal wave functions of Slepian; cf. [CL2, §2.2] and [Sle, SP].

The geometric harmonics now allow one to define the extension of $f : X_A \rightarrow \mathbb{R}$ by

$$\hat{f}(\mathbf{x}_i) := \sum_{\ell=1}^{|A|} \langle f, \psi_\ell \rangle_A \Psi_\ell(i), \quad (4)$$

where $\langle f, g \rangle_A = \sum_{i \in A} f(i)g(i) = f^T g$ is the restriction of the usual inner product to A . Note that (4) makes sense for *any* row \mathbf{x}_i of X , and in the case when $i \in A$, formula (4) simply recovers the representation of the vector f in terms of the ONB $(\psi_\ell)_{\ell=1}^{|A|}$.

2.2. The iteration scheme. Consider the j^{th} coordinate function $\gamma_j(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{e}_j = x_{ij}$ (where \mathbf{e}_j is the standard unit basis vector with 1 in the j^{th} entry and 0 elsewhere), and let A_j consist of those rows of X for which x_{ij} is *not* a missing value. To address the problem of missing data, γ_j can be modeled as a function of the characteristics $m = 1, \dots, j-1, j+1, \dots, p$. In this way, we can think of the j^{th} column of X as a partially-defined function $\gamma_j : X_{A_j} \rightarrow \mathbb{R}$ to which the extension algorithm of geometric harmonics may be applied. This heuristic appears to be valid in practical examples, especially in applications to natural datasets with high dimensionality; cf. §3.

The iteration scheme is initialized by stochastically imputing the missing values in the dataset. More precisely, for each j , the missing values of γ_j are drawn from a normal distribution with mean and variance computed equal to the sample mean and sample variance of $\gamma_j|_{A_j}$. For each iteration, the following steps are conducted for each $j = 1$ to d :

- (1) Form a new matrix $X^{(j)}$ by deleting column j from X .
- (2) Compute the (restricted) Gram matrix $K^{(j)}$ by applying (1) to the rows of $X^{(j)}$.
- (3) Define A_j to be the set of rows of X for which γ_j is defined (i.e., for which x_{ij} is *not* a missing value).
- (4) Compute the eigendata of $K^{(j)*} K^{(j)}$, i.e., find $\lambda_\ell^{(j)} \in \mathbb{R}$ and $\psi_\ell^{(j)} \in \mathbb{R}^{|A_j|}$, that satisfy

$$\sum_{m \in A_j} K(i, m) \psi_\ell^{(j)}(m) = \lambda_\ell^{(j)} \psi_\ell^{(j)}(i),$$

for $\ell = 1, \dots, |A_j|$. Note that these functions $\psi_\ell^{(j)}$ are defined only for $i \in A_j$.

- (5) As in (3), construct the *geometric harmonics* $\Psi_\ell^{(j)}$ by

$$\Psi_\ell^{(j)}(i) := \frac{1}{\lambda_\ell^{(j)}} \sum_{m \in A_j} K(i, m) \psi_\ell^{(j)}(m).$$

Note that these functions $\Psi_\ell^{(j)}$ are defined for all i .

(6) As in (4), fill in the missing values of γ_j using

$$\hat{\gamma}_j(i) = \sum_{\ell=1}^{|A_j|} \langle \gamma_j, \psi_\ell^{(j)} \rangle_{A_j} \Psi_\ell^{(j)}(i).$$

Remark 1 (Random shuffle of characteristics). In an implementation of this method, the order in which the characteristics are considered is permuted after each iteration. In other words, let σ be a random permutation of the characteristics $\{1, \dots, p\}$. Then, in each of steps (1)–(6), the index j is replaced by the random index $\sigma(j)$. This step prevents introducing a bias in the degree of correction corresponding to different indices. More precisely, we found during experimentation that introducing this step caused the algorithm to require 1–2 more iterations to achieve optimal convergence, but that final results were more accurate (in terms of the measurements discussed in section §3). We omitted this detail in the description above, in an attempt to keep the notation from becoming overly heavy.

Remark 2 (Numerical difficulties). It is clear from (3) that the extension procedure can be ill-conditioned for $\lambda_\ell \approx 0$; see [Laf, §3.2]. This is taken into account in the implementation of the method.

3. DISCUSSION AND ANALYSIS OF RESULTS

3.1. Synthetic data. To simulate data with intrinsic nonlinear geometry, we generated 250 points on a swiss roll in \mathbb{R}^3 . The points along the spirals were parameterized by arc length to ensure even spacing:

$$ds = \frac{1}{2} \left(\sqrt{1+t^2} + \sinh^{-1}(t) \right) dt.$$

We added height by generating 5 of these spirals on top one another, also evenly spaced. To enrich the dimensionality, the swiss roll was embedded into \mathbb{R}^{30} , rotated in many random directions, and then Gaussian noise was added to ensure the dataset did not lie in any linear subspace \mathbb{R}^k with $k < 30$. The spread (rate of increase of distance to spiral), height, noise (variance), and number of rotations were chosen experimentally. This test dataset was designed to have such a uniform distribution of points so as to avoid any confounding influence of anisotropy while testing the method.

Figure 1 shows the results of some experiments with a synthetic dataset $X = (\mathbf{x}_i)_{i=1}^{250} \subseteq \mathbb{R}^{30}$. The dataset X was run through a program which deletes any value with fixed probability p . To assess the effectiveness visually, we used diffusion mapping $D_t : \mathbb{R}^{30} \rightarrow \mathbb{R}^3$ to embed the output of our experiments and produce 3-dimensional plots. For information on diffusion mapping (also introduced in Lafon’s thesis; cf. [Laf, §2]), see the enjoyable introduction [LL] and further references [CL1, KCLZ, NLCK1, CKL⁺, NLCK2]. We used a modified version of the excellent implementation of diffusion mapping provided by Laurens van der Maaten in his “Matlab Toolbox for Dimensionality Reduction” which can be found at

<http://homepage.tudelft.nl/19j49/>.

The imputed datasets in Figure 1 visually match the original datasets quite well (although this is much easier to see while rotating the figures on the computer).

Let $\mathbf{x}_i^{(k)}$ denote the imputed version of the i^{th} datapoint \mathbf{x}_i at iteration k of the algorithm described in §2.2. In order to track the effectiveness of the iterated geometric harmonics method, several points with missing data were selected at random. For each such datapoint, we drew line segments connecting $\mathbf{x}_i^{(0)}$ (plotted as \star) to both $\mathbf{x}_i^{(end)}$ (plotted as \bullet) and \mathbf{x}_i (plotted as \blacksquare). Here, $\mathbf{x}_i^{(0)}$ means the stochastically initialized version of the point computed before beginning steps (1)–(6), and $\mathbf{x}_i^{(end)} = \mathbf{x}_i^{(10)}$ is the final imputed version of the point (the experiment was run for 10 iterations). This provides a graphical indication of how far the algorithm moves some points in order to restore them to their (almost) correct position. In most cases, the imputed point $\mathbf{x}_i^{(end)}$ is so close to the original point \mathbf{x}_i that the two

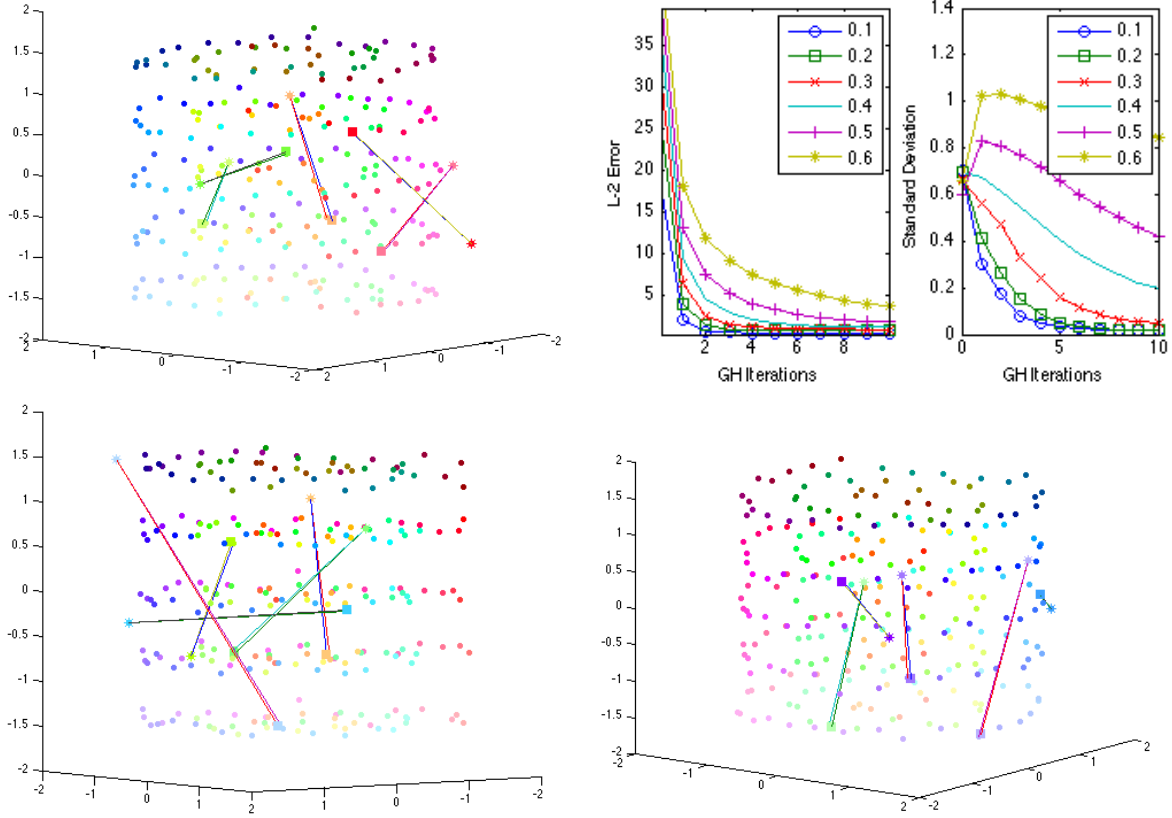


FIGURE 1. The three scatter plots show the swiss roll dataset, reconstructed after annihilating data at rate $p = 0.4$; see §3.1. The final imputed versions of each data point is plotted as a dot \bullet . However, for each plot, 5 data points \mathbf{x}_i , $i = 1, \dots, 5$, were selected for observation at random. The stochastically initialized point $\mathbf{x}_i^{(0)}$ is plotted as a \star , and the final imputed point $\mathbf{x}_i^{(end)}$ is plotted as a \blacksquare . The graphs at top right show the decay of L^2 error (left; lower values are better) and standard deviation (right) with each iteration; see §3.1.

are difficult to distinguish. It is clear from the figures that points often required a large correction, and that the algorithm was able to supply this correction.

Since we have the original intact dataset to compare with, we were able to compute the discrepancy between the original and imputed images using dimension-scaled standard L^2 -norm

$$\text{Error} = p^{-1/2} \sqrt{\sum_{i,j} \left(X(i,j) - \tilde{X}(i,j) \right)^2}, \quad (5)$$

where X is the original dataset, \tilde{X} is the imputed dataset, and d is the number of dimensions, or parameters, of each data point in the dataset.

The graphs in Figure 1 show the L^2 error and standard deviation for values of $p = 0.1, 0.2, \dots, 0.6$. Even with a data annihilation rate of $p = 0.50$ (in which case approximately 50% of the data is lost), the L^2 error decreases by an order of magnitude in about 4 or 5 iterations, after which it decreases only slightly.

3.2. Natural data.

3.2.1. Image datasets. The method was tested on both the Olivetti and UMIST faces datasets, with excellent results. The Olivetti dataset is comprised of 400 photos (10 shots of each of 20 people), each of which is 64×63 pixels. The UMIST dataset is comprised of 15–30 photos of each of 20 people, each of which is 112×92 pixels. Both datasets are freely available on the web.

To test the method, the test dataset was run through a program which deletes any value in the matrix with probability p . For example, Figure 2 shows the result of running this method on the entire Olivetti faces dataset, with data annihilation rate $p = 0.70$. Each pixel has probability $p = 0.70$ of being deleted, so approximately 70% of the data is lost. Nonetheless, as shown in Figure 3, IGH is able to reconstruct the data with a startlingly high degree of accuracy: while some noise remains (inevitably), the people in the photos are now clearly recognizable.

The same procedure was carried out for the UMIST faces database. Figure 5 shows the original, damaged, and imputed (reconstructed) versions of one of the faces in the dataset.

Remark 3. The image reconstruction scheme presented here is novel in that a pixel’s value is imputed *without making any direct comparison to neighboring pixel values within the same image*. Instead, the pixel’s value is inferred based entirely on comparisons with *other* images. In other words, the algorithm makes no use of the fact that the pixel in position (12, 37) in image i is adjacent to pixels (11, 37), (13, 37), (12, 36), or (12, 38) of image i , and that consequently, these pixels values are likely correlated.

To study the effect of sample density on the performance of the method, we used images from the COIL-20 dataset of 128×128 pixel photos of 20 different objects. The full COIL-20 dataset contains 72 photographs of each object, corresponding to increments of 5° of rotation. For the experiment reported in Figure 6, each test was performed using 8 images of the same object: a toy car (object #3). For “sparsity 1”, we used a sequence of successive images (5° rotation between samples), for “sparsity 2”, we used a sequence of every-second images (10° rotation between samples), and so forth. The graph indicates that error increases sharply as sparsity varies from 1 to 3, after which point sparsity makes little difference.

3.2.2. Weather data. The algorithm was tested on a set of 2000 data points recorded at San Diego Lindbergh Field’s weather station GHCND:USW00023188 (located at the San Diego International Airport: latitude 32.73361, longitude -117.18306) during the period Jan. 1–Mar. 24 of 2010. Each record contains 25 parameters, including average and quantile measurements of temperature, pressure, dew point, wind velocity, and cloud cover, and these measurements were recorded hourly (24 times per day). As with the image datasets, individual measurements were deleted with probability p , and the resulting decimated dataset was used as input for the IGH algorithm. Results are presented in Figure 7 and Figure 8; in 2–3 iterations, the algorithm reduces the relative error of a random initialization by approximately 70%. Note that the random initialization is stochastically imputed by drawing from a normal distribution with mean and variance determined by the nonmissing data, exactly as is done during multiple imputation (MI) routines. The standard deviation is not as well-behaved for this experiment, but note that the scaling of the axes magnifies what are actually very small discrepancies.

3.3. Convergence rates. The convergence properties of the IGH method are difficult to study analytically, as can be seen from representing the updates in terms of matrix multiplication, as follows. Let the j^{th} column of X after t iterations be denoted by $\mathbf{x}_j^{(t)}$. A single iteration of IGH is comprised of doing each of the updates

$$\mathbf{x}_j^{(t+1)} \leftarrow L_j^{(t)} \mathbf{x}_j^{(t)}, \quad (6)$$

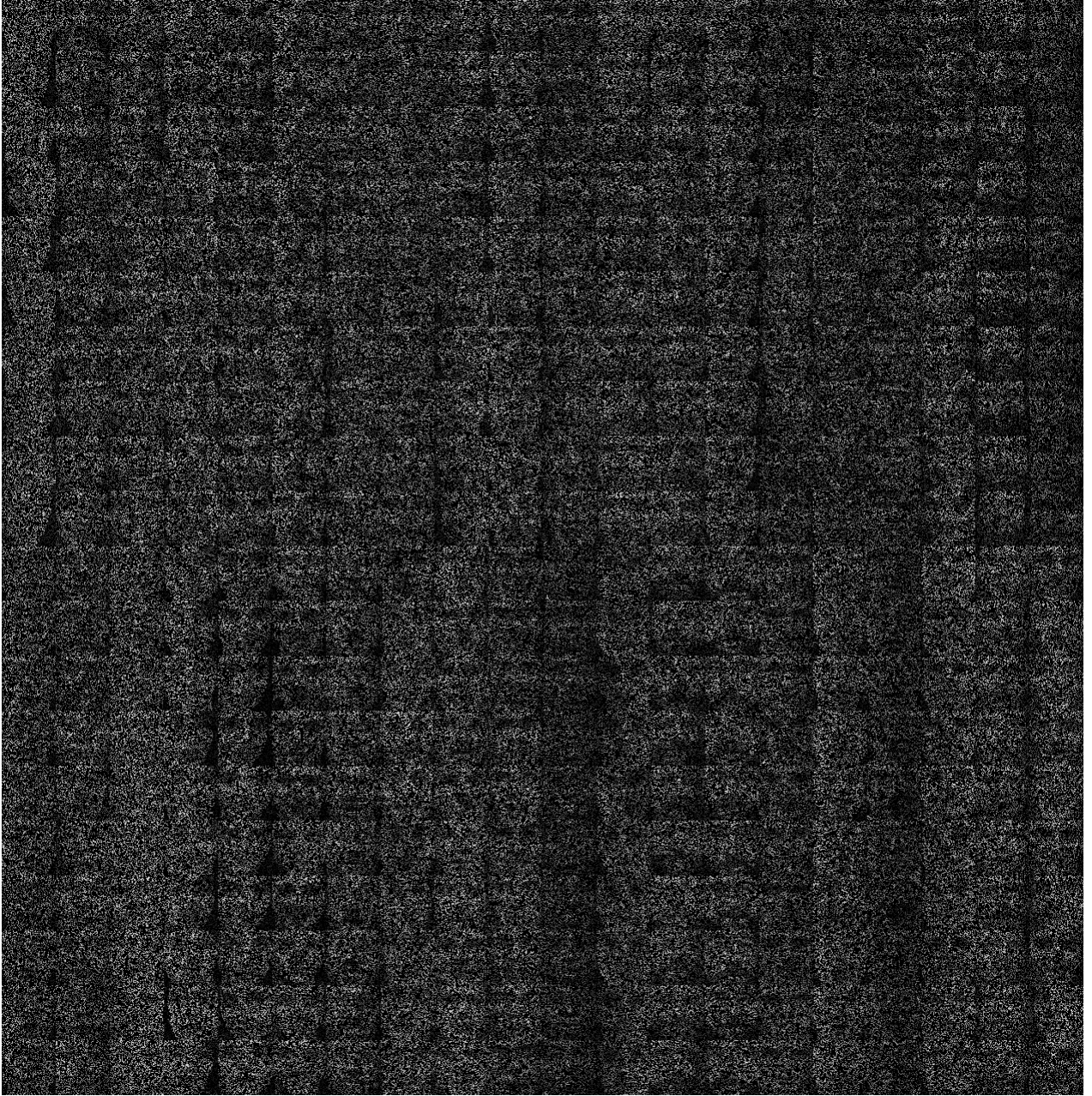


FIGURE 2. The Olivetti faces dataset with data annihilation rate of $p = 0.70$; see §3.2. Missing data values have been plotted as black pixels, and this is the reason for the apparent poor image quality. This is the data used as input for the method; the resulting output is displayed in Figure 3.

in some random order (see Remark 1), where the linear operators $L_j^{(t)} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ update the j^{th} column on iteration t according to

$$L_j^{(t)} = \Psi^{(j,t)}(\psi^{(j,t)})^T K_{A_j}^*, \quad j = 1, \dots, |A_j|.$$

Here, $\Psi^{(j,t)}$ is the version of $\Psi^{(j)}$ appearing at time t , where $\Psi^{(j)}$ is the $n \times |A_j|$ matrix whose columns are the geometric harmonics $\psi_\ell^{(j)}$. It is shown in [Laf, CL2] that for $A \subseteq X$, geometric harmonics produces the extension $F : X \rightarrow \mathbb{R}$ of $f : A \rightarrow \mathbb{R}$ which is maximally

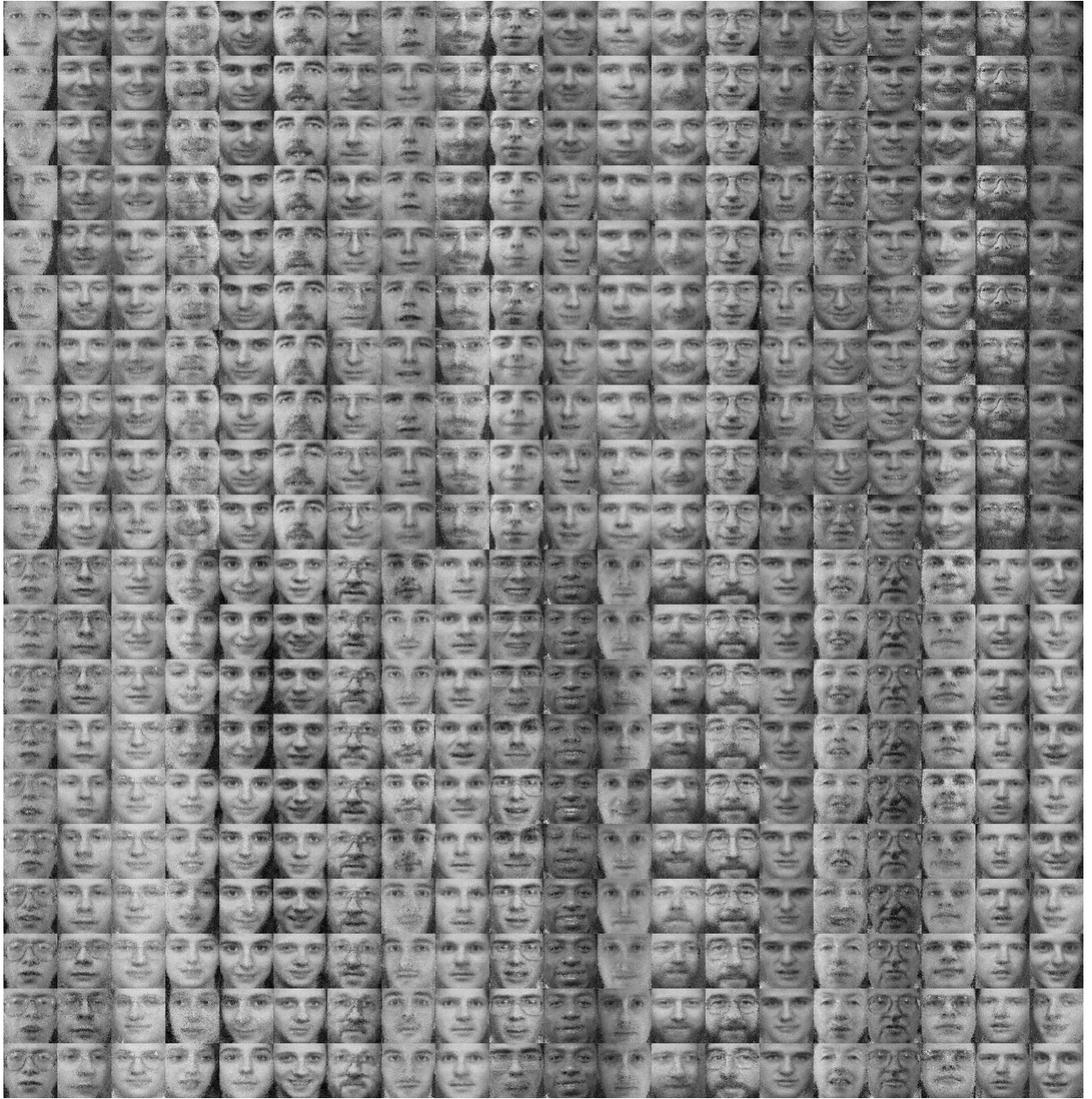


FIGURE 3. The dataset after imputing missing values using the technique of iterated geometric harmonics; see §3.2. This is the output that resulted from using the data shown in Figure 2 as input for the method.



FIGURE 4. For comparison, this figure shows the original (undamaged) photos from the bottom row of the Olivetti dataset displayed in Figure 3.

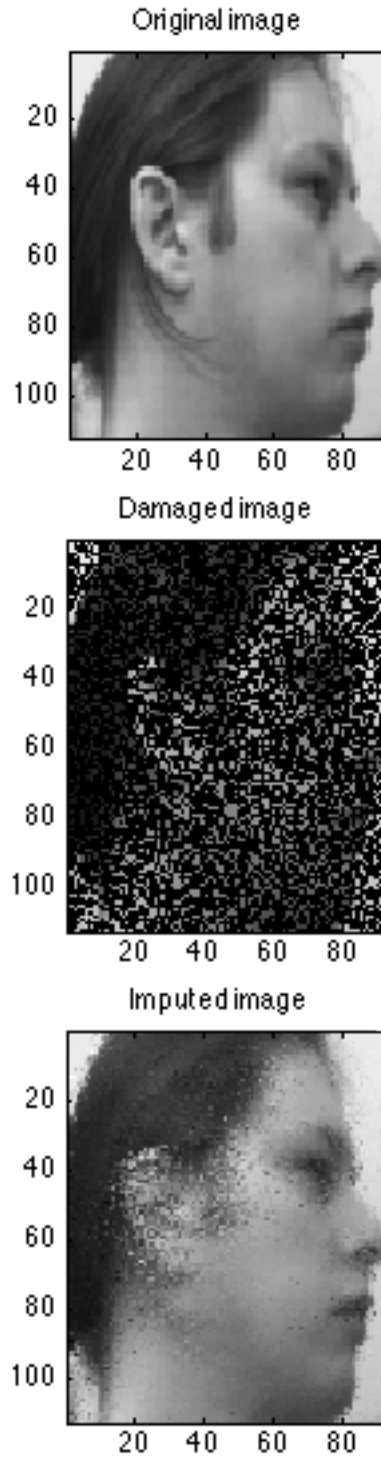


FIGURE 5. The original, damaged, and imputed (reconstructed) versions of one of the faces in the UMIST dataset. These images are 112×92 pixels and there were 30 images in this cluster (the dataset contains 30 photos of this particular individual). This image highlights two areas where the method has difficulty: (i) divisions between very dark and very light regions, and (ii) areas with larger movement (because of camera angle, the ear moves faster across the frame than nose and lips).

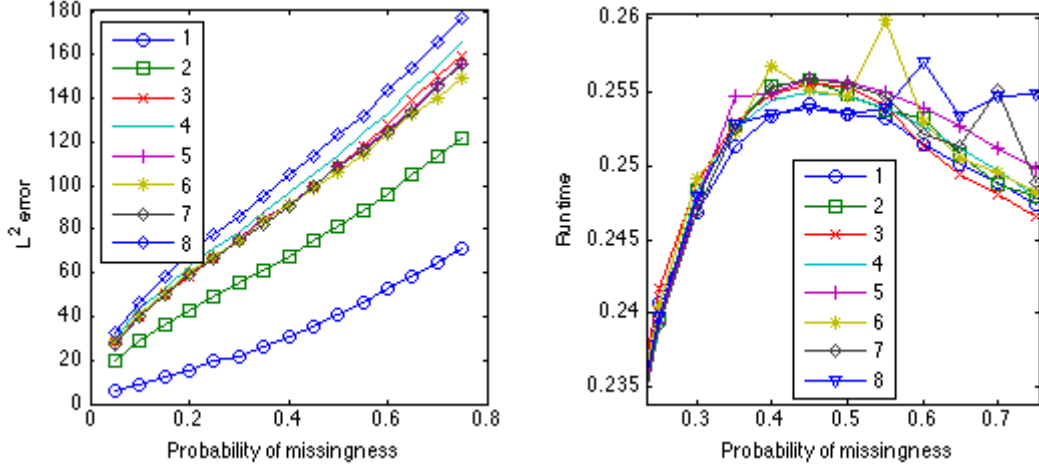


FIGURE 6. Performance for the COIL-20 dataset. Left: Dependence of accuracy on data sparsity (lower values are better); see §3.2.1. Right: Dependence of run time on data sparsity (lower values are better); see §4.1. Run time continues to drop sharply as p tends to 0.

concentrated on A in the sense that the energy of F on $X \setminus A$ is minimized. Consequently, each update (6) replaces $\mathbf{x}_j^{(t)}$ with a new version $\mathbf{x}_j^{(t+1)}$ that has minimal energy on $X \setminus A_j$. This suggests formulation of a condition for convergence based on a norm which incorporates this energy contribution; every $L_j^{(t)}$ is nonexpansive for

$$\|\mathbf{x}_j\| := \|\mathbf{x}_j\|_{L^2(X_{A_j})} + \mathcal{E}_{X \setminus X_{A_j}}(\mathbf{x}_j),$$

where $\mathcal{E}_{X \setminus X_{A_j}}$ is the restriction of the usual graph energy to the complement of A_j .

The rate at which the data converges with geometric harmonics decreases approximately linearly with the rate of missing data (i.e., as $p = 0.1, 0.2, \dots, 0.9$). In general, we found that IGH stabilizes after 2 to 10 iterations, with 4 being sufficient in many cases. The bottom left image of Figure 1 shows the average relative L^2 error versus number of iterations for swiss rolls ranging from 10 to 60 percent missing values, and the right image in the figure shows the variance of each curve. The general trend is as expected: datasets with more missing values require more iterations to stabilize and are less consistent in their outcome.

4. IMPLEMENTATION

MATLAB (R2013a) code for the IGH algorithm, including the examples discussed above, may be found at

<http://www.calpoly.edu/~epearse/IGH.html>.

Much of this code was written for versatility rather than efficiency, so runtimes could likely be improved. All tests were run on a Mac Pro with 3.2 Ghz Quad-Core Intel Xeon processor, 6GB of 1066 MHz DDR3 RAM, running OS X v10.7.5.

4.1. Execution speed. To test how run time depends on the size of the dataset and on p (the rate of missingness), a number of images were retrieved from the UMIST dataset (so $d = 10,304$) and the IGH procedure was applied to this subset, for different values of p . Figure 9 indicates that run time increases approximately linearly with the number of records in the dataset, and that while run time increases with p , the effect of missingness is negligible for $p > 0.2$. For each experiment, data annihilation was performed randomly (with the specified probability p) and this accounts for the irregularity in the (approximately) linearly increasing graphs.

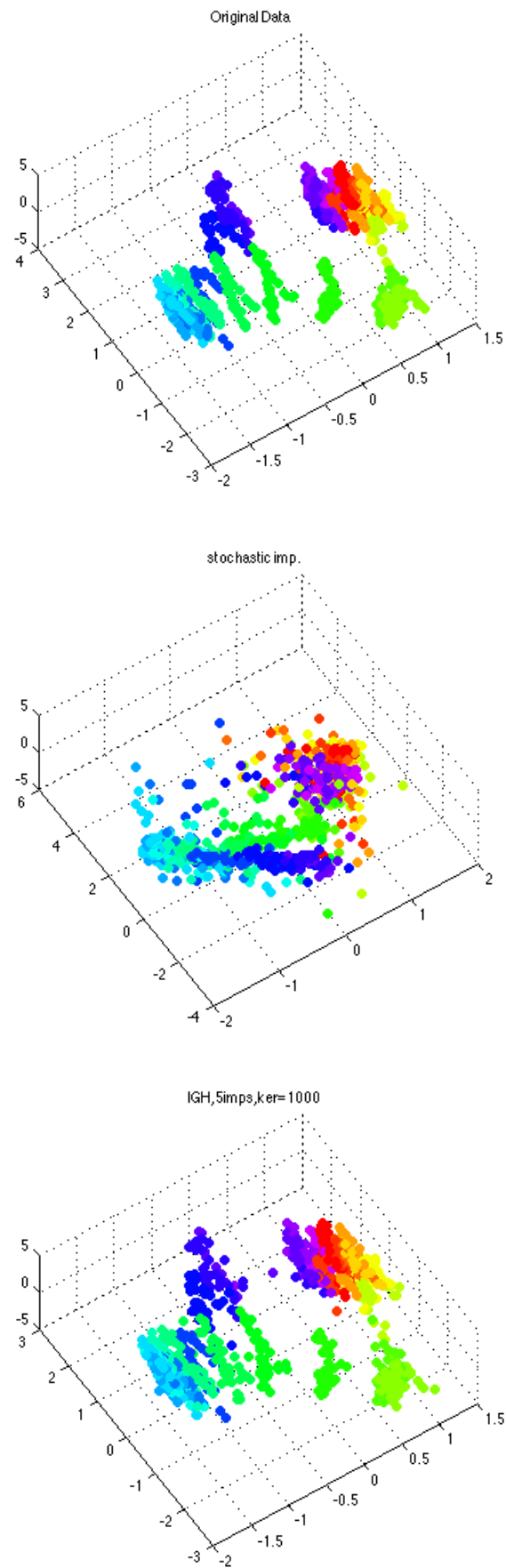


FIGURE 7. The original, damaged, and imputed (reconstructed) versions of the San Diego weather data.

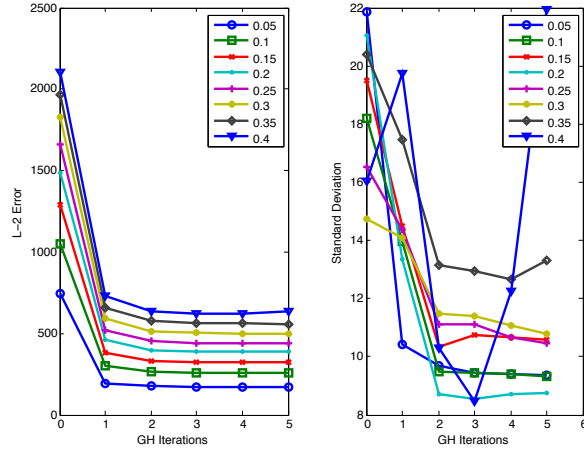


FIGURE 8. Performance of IGH on the San Diego weather data: the decay of L^2 error (left; lower values are better) and standard deviation (right) with each iteration. See §3.2.2.

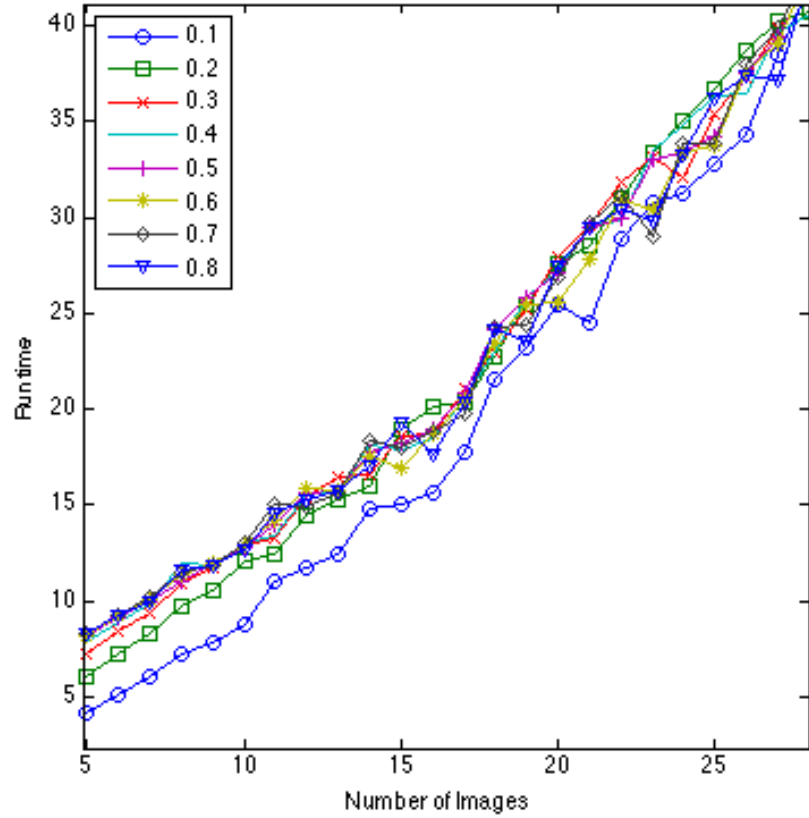


FIGURE 9. Results of the experiment described in §4.1 (lower values are better): run times increase linearly with the number of data points and while slightly shorter for low p -values, change very little with p .

REFERENCES

- [Aro] N. Aronszajn. Theory of reproducing kernels. *Trans. Amer. Math. Soc.* **68**(1950), 337–404.
- [BN] Mikhail Belkin and Partha Niyogi. Laplacian Eigenmaps for Dimensionality Reduction and Data Representation. *Neural Computation* **15**(2003), 1373–1396.
- [BNS] Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: a geometric framework for learning from labeled and unlabeled examples. *J. Mach. Learn. Res.* **7**(2006), 2399–2434.
- [BPV⁺] Y. Bengio, J.F. Paiement, P. Vincent, O. Delalleau, N. Le Roux, and M. Ouimet. Out-of-sample extensions for LLE, Isomap, MDS, eigenmaps, and Spectral Clustering. *Advances in Neural Information Processing Systems* **16**(2006), 214–225.
- [CKL⁺] R. R. Coifman, I. G. Kevrekidis, S. Lafon, M. Maggioni, and B. Nadler. Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems. *Multiscale Model. Simul.* **7**(2008), 842–864.
- [CL1] Ronald R. Coifman and Stéphane Lafon. Diffusion maps. *Appl. Comput. Harmon. Anal.* **21**(2006), 5–30.
- [CL2] Ronald R. Coifman and Stéphane Lafon. Geometric harmonics: a novel tool for multiscale out-of-sample extension of empirical functions. *Appl. Comput. Harmon. Anal.* **21**(2006), 31–52.
- [DM] Petros Drineas and Michael W. Mahoney. On the Nyström Method for Approximating a Gram Matrix for Improved Kernel-Based Learning. *Journal of Machine Learning Research* **6**(2005), 2153–2175.
- [End] Craig K. Enders. *Applied Missing Data Analysis*. Methodology in the Social Sciences. The Guilford Press, 2010.
- [FBCM] Charless Fowlkes, Serge Belongie, Fan Chung, and Jitendra Malik. Spectral grouping using the Nyström method. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **26**(2004), 214–225.
- [FBM] Charless Fowlkes, Serge Belongie, and Jitendra Malik. Efficient spatiotemporal grouping using the Nyström method. *IEEE Comput. Vision Pattern Recogn.* (Dec. 2001).
- [GM] Alex Gittens and Michael W. Mahoney. Revisiting the Nyström Method for Improved Large-Scale Machine Learning. *CoRR* **abs/1303.1849**(2013).
- [KCLZ] Yosi Keller, Ronald R. Coifman, Stéphane Lafon, and Steven W. Zucker. Audio-visual group recognition using diffusion maps. *IEEE Trans. Signal Process.* **58**(2010), 403–413.
- [Laf] S. Lafon. *Diffusion maps and geometric harmonics*. PhD thesis, Yale University, 2004.
- [LL] S. Lafon and A.B. Lee. Diffusion maps and coarse-graining: a unified framework for dimensionality reduction, graph partitioning, and data set parameterization. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **28**(Sept 2006), 1393–1403.
- [LR] Roderick J A Little and Donald B Rubin. *Statistical Analysis with Missing Data*. Wiley-Interscience, 2002.
- [NLCK1] Boaz Nadler, Stéphane Lafon, Ronald Coifman, and Ioannis G. Kevrekidis. Diffusion maps—a probabilistic interpretation for spectral embedding and clustering algorithms. In *Principal manifolds for data visualization and dimension reduction*, volume 58 of *Lect. Notes Comput. Sci. Eng.*, pages 238–260. Springer, Berlin, 2008.
- [NLCK2] Boaz Nadler, Stéphane Lafon, Ronald R. Coifman, and Ioannis G. Kevrekidis. Diffusion maps, spectral clustering and reaction coordinates of dynamical systems. *Appl. Comput. Harmon. Anal.* **21**(2006), 113–127.
- [SSSM] Bernhard Schölkopf, Alexander Smola, Er Smola, and Klaus-Robert Müller. Nonlinear Component Analysis as a Kernel Eigenvalue Problem. *Neural Computation* **10**(1998), 1299–1319.
- [SP] D. Slepian and H. O. Pollak. Prolate spheroidal wave functions, Fourier analysis and uncertainty. I. *Bell System Tech. J.* **40**(1961), 43–63.
- [Sle] David Slepian. Prolate spheroidal wave functions, Fourier analysis and uncertainty. IV. Extensions to many dimensions; generalized prolate spheroidal functions. *Bell System Tech. J.* **43**(1964), 3009–3057.
- [vB] Stef van Buuren. *Applied Missing Data Analysis*. Chapman & Hall/CRC Interdisciplinary Statistics. Chapman and Hall, 2012.
- [WS] C. Williams and M. Seeger. Using the Nyström Method to speed up kernel machines. *Neural Inf. Process. Systems* **13**(2001), 682–688.